

基于 QUEST 决策树的遥感影像土地利用分类 ——以云南省丽江市为例

吴健生^{1,2}, 潘况一³, 彭 建², 黄秀兰¹

(1. 北京大学深圳研究生院城市规划与设计学院, 城市人居环境科学与技术重点实验室, 深圳 518055;

2. 北京大学城市与环境学院地表过程与模拟教育部重点实验室, 北京 100871;

3. 浙江省测绘大队, 杭州 310030)

摘要: 土地利用分类精度直接决定土地利用/土地覆被变化相关研究的准确性, 而基于决策树的遥感影像分类是近年来提高土地利用分类精度的重要方法。QUEST 决策树在影像解译和空间表达方面, 运算速度和分类精度均优于普通 CART 等决策树方法。本文以云南丽江地区为例, 应用 QUEST 决策树分类方法, 对该地区的 Landsat TM 5 影像图进行分类, 同时将地形因素、植被指数作为地学辅助数据的因子添加到分类波段中, 进行不同特征融合, 来处理目标类别间的非线性关系, 该方法在处理图像理解知识方面具有更大的灵活性; 同时与普通决策树分类法的遥感影像分类的结果相比较, Kappa 系数值从原来的 0.789 提高到 0.849。在地形复杂的山地地区, 针对 TM 影像数据, 选择基于 QUEST 决策树分类能够有效提高土地利用分类结果精度。

关键词: 土地利用分类; 决策树; 地学辅助数据; 分类精度; 光谱特征

文章编号: 1000-0585(2012)11-1973-08

1 引言

在土地利用类型分类中应用遥感图像数字处理技术, 可以大幅度提高影像解译和空间表达的效率, 为大面积土地覆盖和土地利用专题图的制作提供了可能^[1]。其中, 遥感影像分类的精度在很大程度上决定了土地利用/土地覆被变化研究的准确性。影响遥感影像分类精度的因素主要有两个: 一是对混合像元, 即在分类结果中处于不同类别边缘的像元的处理方式; 二是对光谱特征变异, 即影像上地类的光谱特征在环境背景的影响下产生变异的处理, 特别是存在“同物异谱”和“同谱异物”现象时^[2]。

目前, 关于遥感影像分类的研究从方法上可以分为两类。一类是基于像素的遥感影像分类方法, 包括多重滤波、基于主成份分析光谱角度制图的分类模型、人工神经网络(ANN)、决策树分类等^[3]。此类方法主要是针对 TM 和 SPOT 等多光谱遥感影像, 尽管经过多年的研究改进, 但在分类结果中仍存在“椒盐效应”, 并且需要大量分类后处理工作进行结果修正^[4]。另一类方法采用面向对象的研究方法, 在一定的同质性标准下, 通过影像分割获得基元, 利用影像的纹理、邻域信息、GIS 辅助数据在模糊分类思想的指导下确定分割对象的所属类别, 此种方法多应用于高分辨率遥感影像, 而在山地地区大范围资

收稿日期: 2011-12-12; 修订日期: 2012-04-23

基金项目: 国家科技支撑计划资助项目 (2008BAB38B03)

作者简介: 吴健生 (1965-), 男, 副教授, 湖南人, 研究领域为景观生态与 GIS。E-mail: wuj@sypku.edu.cn

源遥感调查的应用中, 研究结果精度尚需进一步提高。

以上两种方法的共同趋势是越来越多地依靠多维遥感信息复合来提高遥感影像分类的精度。通过决策树分类法加深对多维遥感信息的认识, 可以更充分的发掘与利用遥感影像数据中隐藏的丰富知识^[5]。决策树是遥感图像分类中的一种分层处理结构, 通过一些判断条件对原始数据集逐步进行二分和细化, 并通过地学辅助数据以及训练样区自身的光谱结构特性, 以一种知识驱动分类的方式完成遥感影像的分类^[6]。决策树具有非参数的特点, 能够处理噪声数据、辅助自动选取特征。在每一级树的划分过程中, 可以使漏分误差和错分误差最小化。决策树分类在许多关于影像分类的行业部门中得到应用, 例如土地覆被制图、森林遥感调查、森林病虫害区划、生物栖息地分类以及土壤景观制图等, 但如何利用地学辅助信息和影像训练区自身的特点, 提高复杂地形中不同地类在影像上的区分度, 这方面的研究还较少^[7, 8]。基于上述研究进展, 本文采用基于 QUEST 决策树的影像分类方法, 并与采用普通决策树分类法 (Classification And Regression Tree, CART) 进行精度对比, 以此来检验该方法在地形复杂地区的多光谱遥感影像精度。

2 研究区概况与研究方法

云南省丽江市地处云南省西北部, 北连迪庆藏族自治州, 南接大理白族自治州, 西邻怒江傈僳族自治州, 东与四川凉山彝族自治州和攀枝花市接壤。总面积 2.06 万 km²。丽江市属滇西北金沙江高山峡谷地貌类型, 地势西北高东南低, 金沙江及其支系深切于崇山峻岭之间, 地貌类型复杂破碎^[9]。由于全球气候变暖、丽江旅游引起该地土地利用结构发生变化, 以及该地具有的不同景观生态效应, 本次选择了丽江作为研究区。本文对 2006 年丽江市地区的 Landsat TM 5 影像数据中 7 个波段进行几何校正、影像裁剪后, 在综合考虑土地利用类型分布情况和数据质量的基础上, 选定 3458×1573 像素的区域作为研究区域 (图 1), 研究区内包括的土地利用类型有建设用地、水田、旱地、草地、林地、水域、冰川积雪等。

决策树分类是一种基于空间数据挖掘和知识发现的监督分类方法。通过决策树学习得到分类规则并进行分类, 分类样本属于严格“非参”, 不需要满足正态分布。其优点是结构清晰, 易于理解; 实现简单, 运行速度快, 准确性高; 可以有效地处理大量数据和高维数据^[10]。同时也可以处理非线性关系, 可以充分利用 GIS 数据库中的地学知识辅助分类, 大大提高了分类精度。

本文采用基于快速、无偏、高效统计树 QUEST 算法 (Quick, Unbiased, and Efficient Statistical Tree, QUEST) 的决策树分类方法, 综合 Landsat TM 影像的光谱特征及多源地学辅助数据, 基于相同的实测训练样本, 建立知识驱动的研究区土地覆被分类规

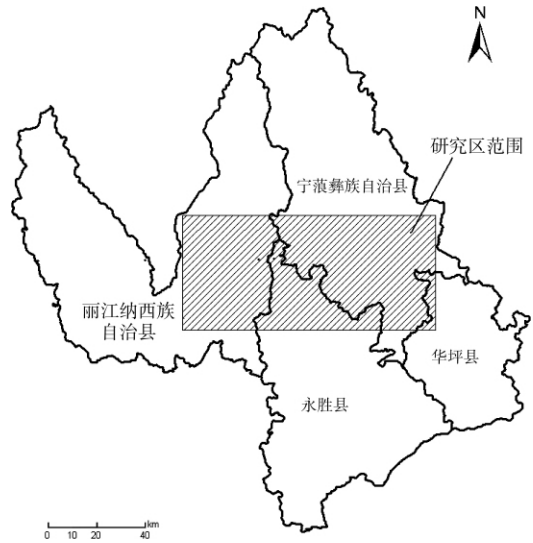


图 1 丽江市研究区范围

Fig 1 The study area of Lijiang

则，并将分类结果与普通决策树分类相比较，建立混淆矩阵进行精度比较，从而实现更为精确的遥感影像分类的方法。普通决策树分类是影像分类数据挖掘算法的一种，将当前数据样本集分为两个子样本集，是一种结构简单的二叉树分类^[11, 12]。

3 影像分类

3.1 训练样本选择

训练样本是进行遥感影像分类的关键，参考 1:10000 比例尺 2004 年丽江市土地利用图和 1:10000 比例尺的数字高程模型图，在影像图上均匀选取各样本训练区（表 1）。

在训练区的土地利用类型分类中，将水田分为两类，原因是一类水田与居民点混淆，混合像元较多，光谱反射率较低，主要位于拉什海的周围，在表中列为零星水田；另一类水田像元较为纯净，光谱反射率较高，在表中列为大片水田。地类名中“其他”代表的是山地地区的阴影部分，参考土地利用图，确认大部分地区覆盖着阴坡植被，本文最后将其归入林地。

3.2 建立基于 QUEST 算法的决策树

3.2.1 QUEST 算法原理 Loh 等在 1997 年提出了 QUEST 决策树构建算法，其基本思想是将变量选择和分割点的选择分开运算，一方面既适用于连续型变量又适用于离散型变量，另一方面考虑到其他一般决策树算法更倾向于选择那些具有更多潜在分割点的预测变量。QUEST 决策树构建在变量选择上基本无偏，同时可通过多个变量构成的超平面在特征空间中区别类别成员和非类别成员。其运算速度和分类精度均优于其他决策树方法^[13~15]。

构建算法流程如下：

(1) 进行预测变量的选择，对所有预测变量 X 与目标变量 Y 进行相关性分析，若 X 为离散变量，使用 χ^2 检验分别计算 X 与 Y 的关联强度，求出 P 值（归入该类的概率）；若 X 是有序或者连续变量，就使用方差分析计算 P 值。

(2) 将所有变量的 P 值与预先设定的阈值， α/M 进行比较 α 为用户指定的显著性水平值在 $(0, 1)$ 之间， M 为预测变量的总数。若均小于阈值，则选择 P 值最小的一个作为分支变量；若均大于阈值，则当 X 为连续或者有序变量时，使用 Levene 方差齐性检验计算 P 值，并在 P 值小于阈值的时刻，选择 P 值最小的一个作为分支变量。如果方差齐性检验的 P 值均大于阈值，则直接选择在第一步中 P 值最小变量作为分支变量。

(3) 若选出的分支变量为离散分类变量，则经过变换，使不同 X 取值时目标变量 Y 取值的差异最大化，并计算其最大判别坐标。

(4) 若 Y 为多分类，则计算 X 的均值，使用聚类分析算法，将这些类别最终合并为两大类，由此将多类类别问题简化为二类判别。

(5) 使用二次判别分析最终明确分割点的位置，并获得所选预测变量 X 的原始取值，最后构建分类规则。

表 1 训练区各地类样本个数及像元数

Tab 1 The number of samples and pixels in the training plots of different ground-objects

地类名	样本个数	样本象元数
河湖水面	23	306
建设用地	36	376
冰川积雪	3	301
林地	45	833
旱地	17	340
零星水田	10	377
大片水田	19	290
草地	5	176
其他	36	376

3.2.2 建立决策树 从选取的训练样本中提取各种光谱和辅助地学特征, 基于上述特征要求, 本文选择 TM 影像的 1~5、7 波段、NDVI 值、DEM 高程值、坡度值, 组成一个 10 波段数据进行分类和精度分析。

TM 的第 6 波段是热红外波段, 一般可以将不放入分类波段, 丽江地区地类与 DEM 高程值和坡度的相关性很高, 因此将其作为重要的辅助地学特征^[16, 17]。水田、旱地、建设用地存在于坡度较低、DEM 值较低的地区, 而冰川积雪、林地、草地存在于坡度较高、DEM 值较高的地区。NDVI 值是反映植被覆盖度的重要指标, 是植物生长状态和植被空间分布密度的最佳指示因子, 并与植被分布密度呈线性相关, 近期被广泛地运用到景观生态、环境监测、农作物估产等领域^[18, 19]。

4 结果分析

4.1 训练区样本分析与决策树建立

本文采用 Jeffries-Matusita (J-M) 距离作为类分离性的度量。训练区样本的在影像中的分离度如表 2。计算的结果表明样本的分离度基本符合分类的要求 (地类之间的分离度大于 1.8 为佳)。

表 2 各地类样本之间 J-M 距离分离度

Tah 2 The J-M separability between two classifications of samples

地类名	分离度	地类名	分离度	地类名	分离度
旱地与草地	1.784	林地与草地	2.000	建设用地与冰川积雪	2.000
建设用地与水田 1	1.848	水田 2 与草地	2.000	其他与冰川积雪	2.000
水田 1 与水田 2	1.947	旱地与水田 2	2.000	冰川积雪与旱地	2.000
建设用地与旱地	1.977	建设用地与水田 2	2.000	冰川积雪与草地	2.000
林地与水田 1	1.979	其他与林地	2.000	河湖水面与水田 2	2.000
旱地与水田 1	1.980	林地与旱地	2.000	其他与水田 2	2.000
林地与水田 2	1.988	建设用地与草地	2.000	冰川积雪与水田 1	2.000
河湖水面与其他	1.991	河湖水面与林地	2.000	冰川积雪与水田 2	2.000
其他与建设用地	1.997	其他与水田 1	2.000	冰川积雪与林地	2.000
建设用地与林地	1.997	河湖水面与水田 1	2.000	河湖水面与冰川积雪	2.000
河湖水面与建设用地	1.997	其他与旱地	2.000	河湖水面与草地	2.000
水田 1 与草地	1.998	河湖水面与旱地	2.000	其他与草地	2.000

共有 3375 个象元样本作为测试变量和目标变量, 利用 ENVI4.7, 建立 QUEST 决策树, 在训练过程中, 决策树的生长深度为 16 层, 结点数为 169 个, 比较方便地产生 IF-Then 形式的规则^[20]。在决策树中依据丽江地区的实际情况以及野外调查数据, 适当修正决策树的划分条件, 并依此为分类方法对多维影像进行分类, 部分区域分类结果如图所示 (图 2)。产生的决策树结构较为复杂, 这里只给出分类生成的决策树结构图 (图 3)。

4.2 精度检验

参考研究区的土地利用数据以及遥感影像图、地形图、居民点分布图、野外调查资料等, 基于研究区地物类型分布面积比例大小, 确定了分层 (Stratified) 随机采样的 2451 个像元样本的地物类型, 分别验证基于普通决策树分类的结果和基于 QUEST 决策树分类的结果。QUEST 决策树分类方法的总精度 90.086%, 比普通决策树分类方法 (85.965%) 高 4.121%; Kappa 系数为 0.849, 比普通决策树分类方法 (0.789) 高 0.060。

其中基于 QUEST 决策树分类方法中的林地、水田、草地的制图精度和用户精度值都得到了显著提高（表 3、表 4）。在使用普通决策树分类方法进行影像分类的过程中，旱地的漏分误差相对较高，而基于 QUEST 决策树进行分类的方法中，水田、旱地的错分率明显减少，是由于在丽江的旱地多位于高海拔地区，地形因素在决策树划分地类中的作用得到明显体现。冰川积雪的光谱反射率很高，并处于高海拔地区，易于与地物类别分离，水体的光谱反射率较低，同样也易于与周围地物区分，因此基于

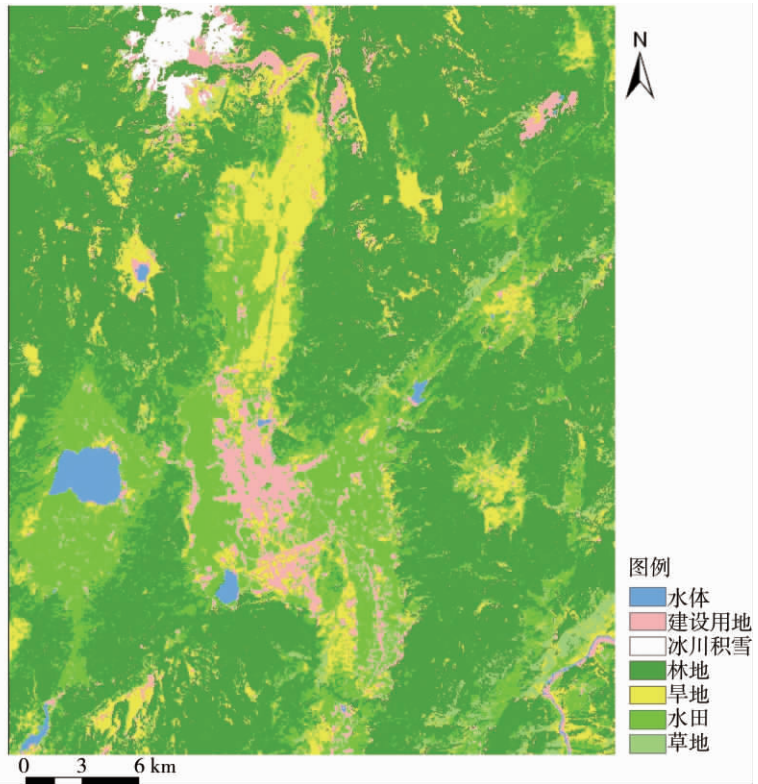


图 2 研究区部分区域的分类结果图

Fig. 2 The classification of part of the study area

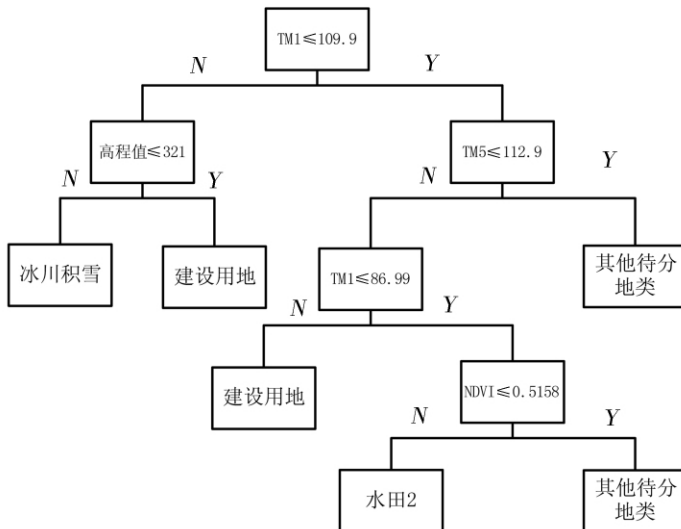


图 3 部分决策树结构

Fig. 3 Part of decision tree

两种方法下的这两类地物的分类精度都很高并且差别不大。林地的分类精度提高幅度不大(制图精度提高了 3.55%),但是林地占该地区的总面积比例最高,因此对总体精度和一致性指数的提高贡献较大。两种方法下的草地错分误差率和漏分误差率均较高,这是由于丽江古城区附近草地与旱地交错分布易混淆,需要通过参考更为精确的辅助数据将其划分,或者通过分解混淆像元的方法将其分离。其混淆矩阵由表 3、表 4 所示,横向地物类别表示实际地物类别像元数,纵向地物类别表示图上地物类别像元数。

表 3 基于普通决策树分类的精度评价

Tab 3 Accuracy assessment on CART-based classification

类别	水体	林地	建设用地	水田	旱地	草地	冰川积雪	总计	用户精度 (%)
水体	95	1	0	0	0	0	0	96	98.96
林地	0	1258	0	2	0	7	0	1267	99.29
建设用地	0	5	343	9	1	19	0	377	90.98
水田	0	86	27	177	61	47	0	398	44.47
旱地	0	1	8	0	173	59	0	241	71.78
草地	0	2	2	0	7	12	0	23	52.17
冰川积雪	0	0	0	0	0	0	49	49	100.00
总计	95	1353	380	188	242	144	49	2451	
制图精度 (%)	100.00	92.98	90.26	94.15	71.49	8.33	100.00		

表 4 基于 QUEST 决策树分类的精度评价

Tab 4 Accuracy assessment on QUEST-based classification

类别	水体	林地	建设用地	水田	旱地	草地	冰川积雪	总计	用户精度 (%)
水体	94	0	0	0	0	0	0	94	100
林地	0	1306	1	0	0	6	0	1313	99.47
建设用地	1	2	343	1	7	13	4	371	92.45
水田	0	33	27	186	16	18	0	280	66.43
旱地	0	1	9	1	216	89	0	316	68.35
草地	0	11	0	0	3	18	0	32	56.25
冰川积雪	0	0	0	0	0	0	45	45	100
总计	95	1353	380	188	242	144	49	2451	
制图精度 (%)	98.95	96.53	90.26	98.94	89.26	12.5	91.84		

5 结论与讨论

本研究在处理遥感影像知识理解和目标识别方面,应用决策树分类方法在很大程度上提高了影像分类结果的精度^[21, 22]。

本研究在地形条件较为复杂的云南丽江地区,基于 QUEST 的决策树分类方法将 NDVI 和地形因子(坡度、坡向)纳入到分类规则相对于单纯地利用光谱信息的遥感影像分类,利用了更多的地学信息。从分类结果上可以得到,基于 QUEST 决策树的遥感影像分类提高了水田、建设用地的分类精度,从而提高了整个研究区影像的总体精度。

今后研究中进行纹理特征、光谱特征与地学辅助信息相互关系的研究,在多维地学信息中发现新的分类规则,构建决策树专家分类模型,反映了遥感影像分类精度提高的一个

方向,具有良好的应用前景。另外,基于QUEST的决策树分类方法还能够快速有效地利用选定的训练样本,从集成遥感影像中获得较为精确的分类规则。在遥感平台上集成多源地学数据,为土地利用的研究夯实牢固的基础^[23]。

参考文献:

- [1] 李小文,刘素红. 遥感原理与应用. 北京:科学出版社,2008. 104~105.
- [2] 韩鹏,龚健雅,李志林,等. 遥感影像分类中的空间尺度选择方法研究. 遥感学报,2010,14(3): 507~518.
- [3] 黎夏,叶嘉安. 基于神经网络的元胞自动机及模拟复杂土地利用系统. 地理研究,2005,24(1): 19~27.
- [4] 苏伟,李京,陈云浩,等. 基于多尺度影像分割的面向对象城市土地覆被分类研究——以马来西亚吉隆坡市中心区为例. 遥感学报,2007,11(4): 521~530.
- [5] 曹丽琴,李平湘,张良培,等. 基于多地表特征参数的遥感影像分类研究. 遥感技术与应用,2010,25(1): 38~44.
- [6] 赵萍,冯学智,林广发. SPOT卫星影像居民地信息自动提取的决策树方法研究. 遥感学报,2003,7(4): 309~315.
- [7] Vaudour E, Carey V A, Gilliot J M. Digital zoning of South African viticultural terroirs using bootstrapped decision trees on morphometric data and multitemporal SPOT images. Remote Sensing of Environment, 2010, 114: 2940~2950.
- [8] 赵清,郑国强,黄巧华. 基于神经网络模型技术的南京市主城区城市森林遥感调查. 地理研究,2006,25(3): 468~476.
- [9] 彭建,李丹丹,张玉清. 基于GIS和RUSLE的滇西北山区土壤侵蚀空间特征分析——以云南省丽江市为例. 山地学报,2007,25(5): 548~556.
- [10] 申文明,王文杰,罗海江,等. 基于决策树分类技术的遥感影像分类方法研究. 遥感技术与应用,2007,22(3): 333~338.
- [11] 潘琛,杜培军,罗艳. 一种基于植被指数的遥感影像决策树分类方法. 计算机应用,2009,29(3): 777~780.
- [12] 朱江洪,李江风,叶菁. 利用决策树工具的土地利用类型遥感识别方法研究. 武汉大学学报:信息科学版,2011,36(3): 301~305.
- [13] Lim T S, Loh W Y, Shih Y S. A comparison of prediction accuracy, complexity, and training time of thirty-three old and new classification algorithms. Machine Learning, 2000, 40: 1~7.
- [14] Loh W Y, Shih Y S. Split selection methods for classification trees. Statistica Sinica, 1997, 7: 815~840.
- [15] 那晓东,张树清,李晓峰,等. 基于QUEST决策树兼容多源数据的淡水沼泽湿地信息提取. 生态学杂志,2009,28(2): 357~365.
- [16] Yang X, Chapman G A, Young M A, et al. Using Compound Topographic Index to Delineate Soil Landscape Facets from Digital Elevation Models for Comprehensive Coastal Assessment. Paper presented at the MODSIM, 2005 International Congress on Modelling and Simulation.
- [17] 李双成,高伟明,周巧富,等. 基于小波变换的NDVI与地形因子多尺度空间相关分析. 生态学报,2006,26(12): 4198~4203.
- [18] 李杭燕. 时间序列NDVI数据集重建方法研究. 兰州:兰州大学硕士学位论文,2010.
- [19] Kross A, Fernandes R, Seaquist J. The effect of the temporal resolution of NDVI data on season onset dates and trends across Canadian broadleaf forests. Remote Sensing of Environment, 2011, 115: 1564~1575.
- [20] 赵萍,傅云飞,郑刘根,等. 基于分类回归树分析的遥感影像土地利用/覆被分类研究. 遥感学报,2005,9(6): 708~716.
- [21] 陈君颖,田庆久. 高分辨率遥感植被分类研究. 遥感学报,2007,11(2): 221~227.
- [22] 邱向红,王周龙,张明明,等. 基于决策树的蓬莱市土地覆盖信息提取. 山东国土资源,2009,25(11): 52~56.
- [23] 陈百明,张凤荣. 我国土地利用研究的发展态势与重点领域. 地理研究,2011,30(1): 1~9.

Research on the accuracy of TM images land-use classification based on QUEST decision tree: A case study of Lijiang in Yunnan

WU Jian-sheng^{1,2}, PAN Kuang-yi^{2,3}, PENG Jian¹, HUANG Xiu-lan¹

- (1. Key Laboratory for Urban Habitat Environmental Science and Technology, School of Urban Planning and Design, Shenzhen Graduate School of Peking University, Shenzhen 518055, Guangdong, China;
2. Key Laboratory for Earth Surface Processes of Ministry of Education, and College of Urban and Environmental Sciences, Peking University, Beijing 100871, China; 3. Zhejiang Brigade of Surveying and Mapping, Hangzhou 310030, China)

Abstract: The accuracy of research on land use/cover change (LUCC) is determined directly by the accuracy of land use classification derived from aerial and satellite images. In analysis of the factors of accuracy of current remote sensing image classification, some methods were introduced to study new trends of classification modes. Some previous studies showed that the speed and accuracy of QUEST (Quick, Unbiased, and Efficient Statistical Tree) decision tree classification were superior to those of other decision tree classifications.

On the basis of this approach, the research classified the Landsat TM-5 images in Lijiang, Yunnan province. This paper compared the result with that of maximum likelihood image classification. The overall accuracy was 90.086%, which was higher than the overall accuracy (85.965%) of CART (Classification And Regression Tree). Meanwhile, the Kappa efficient was 0.849, which was higher than the Kappa efficient (0.760) of CART. Therefore, it is concluded that in the complex terrain area such as in mountainous regions, the choice of QUEST decision tree classification on TM image would improve the accuracy of land use classification.

This type of classification decision tree can precisely obtain new classification rules from integrated satellite images, land use thematic maps, DEM maps and other field investigation materials. Simultaneously, the method can also help users to find new classification rules in multidimensional information, and to build decision tree classifier models. Furthermore, the methods, including a large number of high-resolution and hyperspectral image data, integrated multi-sensor platform, multi-temporal remote sensing image, the pattern recognition and data mining of spectral and texture features, and auxiliary geographic data, will become a trend.

Key words: land use classification; decision tree; geographic auxiliary data; classification accuracy; spectral feature